

# FEDSAL: ENHANCING FEDERATED GRAPH LEARNING THROUGH SALIENCY AWARE CLIENT CLUSTERING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Graph Neural Networks (GNNs) are essential for analyzing structured data but face significant challenges in federated learning (FL) environments, where non-IID client distributions and structural heterogeneity impede convergence and performance. To address these issues, we introduce Federated Saliency Aggregation Learning (FedSal), the first framework to apply saliency maps in GNN-based FL on graph classification tasks. FedSal replaces full-gradient uploads with compact saliency activations, enabling dynamic clustering of clients via simple thresholds ( $\epsilon_{\text{mean}}, \epsilon_{\text{max}}$ ) and cluster-wise model averaging. We further propose FedSal+, which augments node features with positional and random-walk encodings to inject structural priors without exposing raw graph data. Extensive experiments on thirteen molecular, protein, and social-network benchmarks under extreme non-IID splits show that FedSal and FedSal+ achieve higher accuracy, converge faster, and reduce communication cost compared to state-of-the-art methods. These results demonstrate the SOTA performance of saliency-driven clustering for personalized, robust, and communication-efficient federated graph classification tasks.

## 1 INTRODUCTION

Graph Neural Networks (GNNs) have emerged as powerful tools for modeling graph-structured data, achieving state-of-the-art performance in domains including molecular property prediction, clinical risk stratification, community detection, and traffic forecasting Wu et al. (2020); Kipf & Welling (2016); Hamilton et al. (2017); Chen et al. (2020); Cui et al. (2019). Traditional GNN training typically aggregates all graph data centrally, an approach often infeasible due to privacy constraints, corporate policies, and scalability issues. Federated Learning (FL) addresses these challenges by exchanging model updates rather than raw data, preserving data privacy and improving scalability McMahan et al. (2016); Kairouz et al. (2021); Li et al. (2020).

Integrating FL with GNNs, termed Federated Graph Neural Networks (FedGNNs), enhances graph learning across heterogeneous, decentralized datasets. However, FedGNNs face significant hurdles, notably the non-independent and identically distributed (non-IID) nature of decentralized data. Local graphs differ substantially in node features, connectivity, and labels, leading to unstable learning and degraded performance Zhao et al. (2018). This complexity necessitates robust methods for personalized and efficient federated training.

Standard FL methods, such as FedAvg McMahan et al. (2016), inadequately address these heterogeneities, as they fail to consider structural and semantic divergences Xie et al. (2021). Clustered FL approaches, which group clients by compatible updates, have shown promise in enhancing personalization and reducing cross-client interference Sattler et al. (2020). For example, Graph Clustered Federated Learning (GCFL) Xie et al. (2021) clusters clients based on raw gradient similarities, yet raw gradients remain inherently noisy, sensitive to scaling, and prone to variance.

To overcome these limitations, we propose **FedSal**, a novel FedGNN framework leveraging saliency activation maps—gradients with respect to model outputs rather than data—to cluster clients effectively. Saliency maps highlight influential features contributing to predictions, providing more stable and representative client summaries compared to raw gradients. Although saliency maps have demonstrated efficacy in interpretability, robustness, and improved performance across various graph tasks Pei et al. (2024); Wang et al. (2024), their application within Federated Graph Neural Networks (FedGNNs) has not previously been explored. FedSal dynamically clusters clients based on semantic

054 similarity in saliency profiles, enhancing personalized aggregation, model stability, and learning  
 055 efficiency under extreme non-IID conditions

056 Specifically, FedSal clusters clients by computing compact, normalized saliency summaries and parti-  
 057 tioning a cosine-affinity graph via a Stoer–Wagner minimum-cut algorithm with adaptive thresholds  
 058 ( $\epsilon_{\max}$ ,  $\epsilon_{\text{mean}}$ ). Aggregating parameters within clusters mitigates conflicts arising from structural  
 059 divergences and accommodates client-specific preferences. Extending this approach, **FedSal+** incor-  
 060 porates positional and random-walk encodings into saliency maps, injecting structural priors to refine  
 061 clustering resolution and further improve accuracy.

062 Our research addresses four central questions: **RQ1**: *Do saliency activation profiles better capture*  
 063 *inter-client similarities compared to raw-gradient or spectral methods?* **RQ2**: *Does injecting*  
 064 *positional and random-walk structural encodings into saliency summaries improve cluster coherence*  
 065 *and downstream accuracy?* **RQ3**: *What is the trade-off between communication overhead (saliency*  
 066 *summary size) and convergence performance compared to baselines?* **RQ4**: *How do adaptive*  
 067 *thresholds ( $\epsilon_{\max}$ ,  $\epsilon_{\text{mean}}$ ) influence cluster stability, convergence robustness, and fairness for*  
 068 *minority-label clients?*

069 To address **RQ1** and **RQ4**, FedSal computes normalized saliency summaries each round, partitions  
 070 these via adaptive minimum-cut clustering, and averages parameters within clusters, ensuring robust  
 071 handling of dynamic heterogeneity. Addressing **RQ2**, FedSal+ enhances saliency maps with structural  
 072 priors, leading to finer-grained similarity signals, improved cluster coherence, and higher accuracy  
 073 across tasks. For **RQ3**, we systematically analyze communication-performance trade-offs, comparing  
 074 saliency-based methods against gradient- and spectral-based baselines.

075 Experiments across thirteen benchmarks (seven molecular, three protein, and three social network  
 076 datasets) under extreme non-IID splits demonstrate that FedSal and FedSal+ consistently surpass  
 077 baselines in accuracy, convergence speed, and communication efficiency. Ablation studies identify op-  
 078 timal thresholds and cluster counts, highlighting that clients with minority labels benefit significantly  
 079 from our adaptive protocol.

080 Our contributions include:

- 082 • **Unprecedented Use of Saliency Maps in FL and FedGNNs**: No prior work has employed  
 083 saliency maps or FedGNNs. This novel application of utilizing saliency maps to identify  
 084 and cluster clients with similar data feature importance, ensuring that model updates are  
 085 aggregated more effectively and relevantly. could inspire further research.
- 086 • **Introduction of FedSal Architecture**: We propose FedSal, a novel Federated Graph Neural  
 087 Network (FedGNN) architecture that leverages saliency activation maps for client clustering,  
 088 enhancing personalization and model performance in FL environments.
- 089 • **FedSal+**: Structural-prior augmentation using positional and random-walk encodings to  
 090 enhance clustering quality and downstream task performance.
- 091 • **Superior Performance Over State-of-the-Art FedGNNs**: Robust evaluations on real-  
 092 world datasets demonstrating substantial improvements over state-of-the-art FedGNN meth-  
 093 ods.

## 094 1.1 RELATED WORKS

### 095 1.1.1 FEDERATED LEARNING

096 FL was first introduced by McMahan et al. (2016), enabling collaborative training across multiple  
 097 devices under a central server while preserving data privacy Kairouz et al. (2021); Yang et al. (2019);  
 098 Lyu et al. (2022); Yu et al. (2025); Liu et al. (2025b); Chen et al. (2025); Dai et al. (2025); Liu et al.  
 099 (2025a); Tan et al. (2025); Shaikh & Samet (2025); Mai et al. (2024); Fu et al. (2025); Fang et al.  
 100 (2025); Fu et al. (2024); Gao et al. (2024); Wu et al. (2022); Wang et al. (2022b). One of the most  
 101 prominent and standard settings for FL is the FedAvg algorithm by McMahan et al. (2017). This  
 102 algorithm relies on stochastic gradient descent (SGD) based optimization, which affects convergence  
 103 speed and can lead to unstable learning due to its unbiased estimation and averaging of all model  
 104 parameters during aggregation at the central server Zhao et al. (2018); Li et al. (2020); Karimireddy  
 105 et al. (2020). Numerous works have aimed to improve the performance of FedAvg, addressing issues  
 106  
 107

such as performance in heterogeneous (non-IID) settings Wang et al. (2020); Tan et al. (2022a); Chen et al. (2022a); Zhao et al. (2018); Jeong et al. (2018); Huang et al. (2020), communication speeds Hamer et al. (2020), robustness Wang et al. (2019); Yu et al. (2019); Khaled et al. (2020); Liang et al. (2019); Karimireddy et al. (2020); Li et al. (2020), and generalization ability Hamer et al. (2020). One of the most significant challenges in making FL viable for real-world applications is the non-IID data issue, where clients have heterogeneous labels and feature distributions Luo et al. (2021); Tan et al. (2022b); Chen et al. (2022b). Various approaches have been introduced to address this problem, including clustering techniques, which have demonstrated communication efficiencies through group-level personalization of clients. Methods like model-agnostic meta-learning and personalized FL have shown improved generalizability, reducing the effects of non-IID data Fallah et al. (2020); Chen et al. (2018). A recent trend involves decoupling techniques for better personalization, which have also shown improvements in generalizability T Dinh et al. (2020); Li et al. (2021). However, these methods often come with additional communication overhead, although when paired with clustered FL, they have shown reductions in communication costs.

### 1.1.2 FEDERATED GRAPH NEURAL NETWORKS

FL has matured for image and other Euclidean data, but its adoption for graph data remains nascent. Zhang et al. Zhang et al. (2021a) classify FedGNNs into intra-graph, inter-graph, and graph-structured paradigms, while subsequent taxonomies offer finer distinctions He et al. (2021); Fu et al. (2022); Liu et al. (2024). Intra-graph FedGNNs assign each client a subgraph of a larger network to predict missing nodes Zhang et al. (2021b), infer links Chen et al. (2021), or detect communities Baek et al. (2023), with applications in financial crime detection Suzumura et al. (2019). Graph-structured FedGNNs exploit explicit client relationships—e.g., in personalized image processing Chen et al. (2022c) or traffic forecasting Meng et al. (2021). Inter-graph FedGNNs, our focus, train local GNNs on client-specific graph datasets to boost generalization Xie et al. (2021); Zhu et al. (2022); Jiang et al. (2022); Lou et al. (2021).

Recent work has fused core FL techniques and personalization strategies APPLE Luo & Wu (2022), FedCP Zhang et al. (2023), FedGKD Yao et al. (2024), cluster-aware GCFL Xie et al. (2021), gradient-guided FGSSL Huang et al. (2024), and cross-domain FedSSP Tan et al. (2024)—to improve robustness under severe non-IID graph distributions. Complementary advances include two-channel local training Zheng et al. (2021), dynamic client selection Gu et al. (2023), split-channel learning Tan et al. (2023), latent-link generation Xie et al. (2023), embedding masks for personalization Baek et al. (2023), custom optimization losses Guo et al. (2023), meta-learned hyperparameters Wang et al. (2022a), and hierarchical update clustering Briggs et al. (2020). FPL Huang et al. (2023) introduces cluster-aware prototypes to counter domain shifts, while Wang et al.’s FedSLS Wang et al. (2024) aggregates client updates in a saliency–latent space for image-based FL. In the graph domain, centralized saliency-aware regularization has improved robustness Pei et al. (2024), but no prior work leverages saliency maps for client clustering in a federated setting. Our method fills this gap, demonstrating that saliency-guided similarity can effectively drive clustering and aggregation when the data are graphs rather than images.

## 2 METHODOLOGY

### 2.1 TECHNICAL DESIGN

GNNs are powerful tools for learning graph representations and have been widely used in graph mining applications. The model parameters and gradients of GNNs can reflect the underlying graph structure and feature information. Thus, in the FedSal architecture, GNNs are used as the core model for graph mining within the FL framework.

#### 2.1.1 DYNAMIC CLUSTERING WITH SALIENCY AGGREGATION

The FedSal framework introduces a novel approach to dynamically cluster clients by leveraging their transmitted saliency maps. We denote by  $S_i(t)$  the aggregated saliency map of client  $i$  at round  $t$ , and define the update  $\Delta S_i = S_i(t) - S_i(t-1)$ . This mechanism aims to maximize collaboration among homogeneous clients and mitigate the negative effects of heterogeneous clients. When client data distributions are highly heterogeneous, a general FL approach may fail to jointly optimize all local

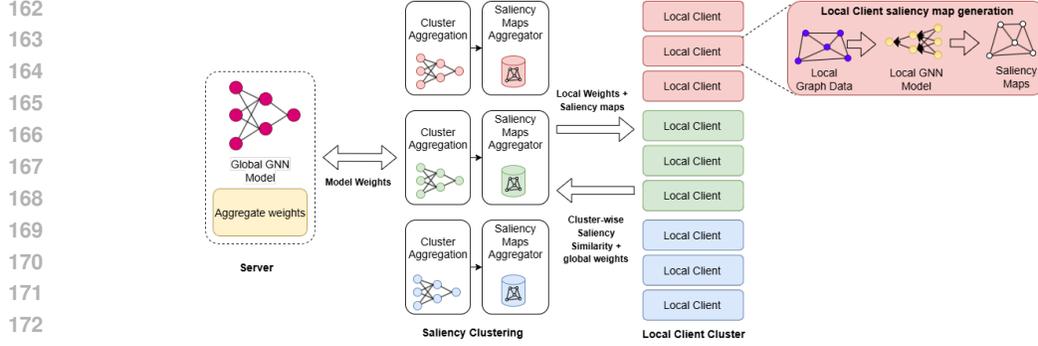


Figure 1: Architectural illustration of FedSal

loss functions. In such cases, after several communication rounds, the general FL algorithm tends to approach a stationary point, indicated by small  $\|\Delta S_i\|$  norms. Therefore, clustering is necessary as the FL process nears that point.

To achieve this, we introduce two hyperparameters,  $\epsilon_{\text{mean}}$  and  $\epsilon_{\text{max}}$ , to guide the clustering process:

**Stopping Criterion for General FL:** The mean norm of saliency-map updates across clients is

$$\frac{1}{N} \sum_{i=1}^N \|\Delta S_i\| < \epsilon_{\text{mean}},$$

where  $N$  is the total number of clients. When this condition holds, the model is considered to be near a stationary point, and further global updates may yield diminishing returns.

**Clustering Criterion:** The maximum norm of saliency-map updates identifies significant heterogeneity:

$$\max_{i=1, \dots, N} \|\Delta S_i\| > \epsilon_{\text{max}} > 0.$$

A high value indicates at least one client whose data distribution deviates substantially, necessitating client clustering to handle diverse data more effectively.

### 2.1.2 MECHANISM OF CLUSTERING AND AGGREGATION

The FedSal framework employs a top-down bi-partitioning mechanism. At each communication round  $t$ , the server receives saliency updates  $\{\Delta S_i\}$  from clients within an existing cluster  $C_k$ . If both the mean and maximum update norms exceed the thresholds  $\epsilon_{\text{mean}}$  and  $\epsilon_{\text{max}}$ , the server constructs a cluster-wise cosine similarity matrix  $\alpha_k$ . Each entry

$$\alpha_{ij} = \frac{\Delta S_i \cdot \Delta S_j}{\|\Delta S_i\| \|\Delta S_j\|}$$

serves as the weight of an edge in a fully connected graph whose nodes represent clients in  $C_k$ . The Stoer–Wagner minimum-cut algorithm is then applied to partition  $C_k$  into subclusters  $\{C_{k1}, C_{k2}\}$ . Further technical details are in Appendix Section 2.1.2.

### 2.1.3 MODEL UPDATE AND AGGREGATION

**Local Update:** Each client  $i$  updates its local model  $\theta_i(t)$  by minimizing the loss on its local data  $D_i$ :

$$\theta_i(t) = \arg \min_{\theta} \text{Loss}(\theta; D_i).$$

After the update, the client computes the per-sample saliency

$$S(x) = \left| \frac{\partial \text{Loss}(f(x; \theta), y)}{\partial x} \right|,$$

and aggregates:

$$S_i(t) = \frac{1}{|D_i|} \sum_{x \in D_i} S(x).$$

216 **Global Aggregation:** For each cluster  $k$ , the server aggregates client updates:

$$217 \theta_{t+1,k} = \theta_{t,k} + \sum_{i \in C_k} \Delta\theta_{t,k,i},$$

219 where

$$220 \Delta\theta_{t,k,i} = \theta_i(t) - \theta_{t,k}$$

221 is the update sent by client  $i$ . Finally, the global model is obtained by averaging across clusters:

$$222 \theta_{t+1} = \frac{1}{K} \sum_{k=1}^K \theta_{t+1,k}.$$

## 227 2.2 THEORETICAL DESIGN

228 This section explores the theoretical and practical effectiveness of saliency activation maps in capturing and reducing structural, feature, and task heterogeneity in graph data using GNNs within a clustered FL setting. To substantiate this, we analyze two problem statements and prove them in Appendix Section A.7 through propositions that saliency maps can effectively represent structural, feature, and task information in the model.

234 **Definition 1 Saliency Map Distortion in FedSal:** Let  $f : \mathcal{G} \rightarrow \mathcal{S}$  be a function mapping from the space of GNN parameters  $(\mathcal{G}, d)$  to the space of saliency maps  $(\mathcal{S}, d')$ . The function  $f$  is considered to have  $\delta$  distortion if for all  $u, v \in \mathcal{G}$ ,

$$235 \frac{1}{\delta} d(u, v) \leq d'(f(u), f(v)) \leq \delta d(u, v).$$

236 This definition ensures that the relationship between graph parameters and their corresponding saliency maps maintains a bounded distortion, thereby preserving the structural and feature-based relationships even after transformation to the saliency space.

240 **Theorem 1 Bourgain’s Embedding Theorem:** Bourgain’s theorem states that any finite metric space  $(X, d)$  can be embedded into Euclidean space with distortion at most  $O(\log n)$ , where  $n$  is the number of points in  $X$ . Bourgain (1985)

244 **Problem 1:** FedSal involves the communication of saliency maps between models with heterogeneous graph structures distributed among different clients. The structure and feature differences can be captured by the saliency maps.

247 **Proposition 1 Structural Sensitivity of Saliency Maps:** Given a model  $M$  with structure represented by the normalized graph Laplacian  $L$ , features  $X$ , and saliency map  $S$ . If we have another model  $M'$  with a different structure  $L'$ , then the saliency-map difference  $\|S' - S\|_2$  is bounded by the structural difference.

251 **Proposition 2 Feature Sensitivity of Saliency Maps:** Given a model  $M$  with structure  $L$ , features  $X$ , and saliency map  $S$ . If we have another model  $M'$  with different features  $X'$ , then the saliency-map difference  $\|S' - S\|_2$  is bounded by the feature difference.

255 **Problem 2:** The communicated saliency maps in FedSal can also capture task heterogeneity.

259 **Proposition 3 Task Sensitivity of Saliency Maps:** Given a model  $M$  with structure  $L$  and features  $X$ . If trained on different tasks, the resulting models will have saliency maps whose differences are bounded.

## 265 3 EXPERIMENTS

### 266 3.1 EXPERIMENTAL CONFIGURATIONS

267 **Datasets and Partitioning:** We utilize 13 graph classification datasets spanning three domains, as detailed in Appendix Section A.5. These include seven molecule datasets (MUTAG, BZR, COX2,

DHFR, PTC\_MR, AIDS, NCI1), three protein datasets (ENZYMES, DD, PROTEINS), and three social network datasets (COLLAB, IMDB-BINARY, IMDB-MULTI). Node features are present in some datasets, and labels are either binary or multi-class in a **graph classification** task.

Two data partitioning settings are used: **Single-dataset setting**: Graphs from a single dataset are randomly distributed among clients, each client receives about 100 graphs, with 10% reserved for testing. **Multi-dataset setting**: Multiple datasets are distributed among clients in groups, with 10% of graphs held for testing.

**Baselines and Model Configuration**: We compare FedSal and FedSal+ with a compact yet representative set of FL baselines, **Selftrain** as the first baseline, the classical aggregators **FedAvg** McMahan et al. (2016) and **FedProx** Li et al. (2018); Yuan & Li (2022), the gradient or cluster-aware methods **GCFL** Xie et al. (2021), the personalised/feature-separation methods **APPLE** Luo & Wu (2022), **FedCP** Zhang et al. (2023), **FedSage** Zhang et al. (2021b), **FGSSL** Huang et al. (2024), Fedstar Tan et al. (2023) and the current cross-domain FedGNN state-of-the-art, **FedSSP** Tan et al. (2024). We omit FedSSP, Fedstar and FedSal+ from IID single-dataset experiments, since their cross-domain cues (domain-adaptive thresholds in FedStar, spectral projections in FedSSP, positional/random-walk encodings in FedSal+) collapse without heterogeneity, and reserve these comparisons for multi-domain benchmarks. The models feature three layers with a hidden size of 64 and are trained using a batch size of 128. An Adam optimizer is employed with a learning rate of 0.001 and a weight decay of  $5 \times 10^{-4}$ .

**Parameter Settings and Computational Resources**: Local epochs are set to 1 for all methods. Clustering hyperparameters  $\epsilon_{\text{mean}}$  and  $\epsilon_{\text{max}}$  are adjusted based on a grid search-based selection for the best performing values. Experiments are conducted on a server with an NVIDIA RTX 4050 GPU, with 16GB of memory, and all experiments are repeated five times for statistical validation.

**Enhancement of FedSal with FedStar**: Inspired by the FedStar framework Tan et al. (2023), we developed **FedSal+**, integrating advanced encoding strategies such as positional and random walk encodings to enhance structural details in FL clients. These enhancements aim to refine client-side feature representations, improving FedSal’s saliency-based model aggregation. By enriching the feature set for each node, FedSal+ produces more detailed saliency maps, focusing on pivotal features during FL, and tests FedSal’s compatibility with methods like structural embeddings. Implementation details are provided in Appendix Section A.4. The primary aim of FedSal+ is to determine if the enriched feature set can enhance FedSal’s effectiveness, particularly in handling non-IID data in multi-dataset federated settings. We hypothesize that these refined feature representations will lead to more precise and informative saliency maps, thereby improving overall model performance in multi-dataset federated environments. This hypothesis sets the stage for our experimental investigations to validate the improved performance metrics that FedSal+ could deliver through advanced structural embeddings.

### 3.2 TEST ACCURACY ANALYSIS

Table 1: Test accuracy (%) results across different FL Methods (mean  $\pm$  standard deviation)

Dataset (# of Clients)	Multi-Dataset			Single-Dataset		
	Molecules (7)	Biochem (8)	Mix (13)	IMDB-BINARY (10)	NCI1 (30)	PROTEINS (10)
Selftrain	75.35 $\pm$ 0.45	70.53 $\pm$ 0.48	69.89 $\pm$ 0.39	76.86 $\pm$ 3.72	62.42 $\pm$ 1.60	73.85 $\pm$ 1.28
FedProx	73.80 $\pm$ 0.65	70.20 $\pm$ 1.10	66.00 $\pm$ 0.90	75.20 $\pm$ 2.90	64.20 $\pm$ 1.60	72.00 $\pm$ 1.20
FedAvg	75.37 $\pm$ 1.21	69.49 $\pm$ 0.58	69.25 $\pm$ 0.77	77.59 $\pm$ 2.45	66.12 $\pm$ 1.33	74.65 $\pm$ 1.18
FedSage	75.90 $\pm$ 0.55	72.40 $\pm$ 0.90	68.00 $\pm$ 0.60	77.80 $\pm$ 2.40	66.50 $\pm$ 1.30	74.80 $\pm$ 1.10
GCFL	76.08 $\pm$ 0.64	70.94 $\pm$ 0.88	69.10 $\pm$ 0.73	78.19 $\pm$ 2.32	65.24 $\pm$ 2.28	75.19 $\pm$ 1.74
APPLE	76.40 $\pm$ 0.70	71.60 $\pm$ 0.75	68.10 $\pm$ 0.80	76.50 $\pm$ 2.60	64.40 $\pm$ 1.30	74.50 $\pm$ 1.15
FGSSL	76.60 $\pm$ 0.60	72.80 $\pm$ 0.70	71.05 $\pm$ 0.50	77.55 $\pm$ 2.20	64.70 $\pm$ 1.40	74.90 $\pm$ 1.10
FedCP	78.40 $\pm$ 0.88	73.10 $\pm$ 0.80	70.95 $\pm$ 0.70	78.00 $\pm$ 2.50	65.00 $\pm$ 1.50	75.00 $\pm$ 1.20
<b>FedSal</b>	76.61 $\pm$ 1.03	<b>73.46 <math>\pm</math> 0.72</b>	<b>71.49 <math>\pm</math> 0.52</b>	<b>80.22 <math>\pm</math> 2.63</b>	<b>69.91 <math>\pm</math> 1.37</b>	<b>75.19 <math>\pm</math> 0.75</b>
FedStar	77.02 $\pm$ 0.39	70.30 $\pm$ 0.61	68.90 $\pm$ 0.89	–	–	–
FedSSP	78.22 $\pm$ 0.70	72.03 $\pm$ 0.65	69.55 $\pm$ 0.55	–	–	–
<b>FedSal+</b>	<b>79.87 <math>\pm</math> 0.86</b>	<b>74.37 <math>\pm</math> 0.44</b>	<b>72.07 <math>\pm</math> 0.61</b>	–	–	–

From the test accuracy results in Table 1. In the single-dataset experiments, FedSal consistently outperforms all competing methods like classical aggregators, personalized approaches and self-training alike, while also exhibiting the most stable performance across clients. In the multi-dataset benchmarks, FedSal+ takes the lead across every setting, clearly beating both general-purpose FL algorithms and state-of-the-art cross-domain techniques such as FedSSP and FedStar. Traditional methods like FedAvg and FedProx struggle under heterogeneous graph distributions, and even specialized solutions (e.g., FedStar, FedCP) cannot match the robustness and scalability delivered by our saliency-based schemes. These trends demonstrate that saliency-driven aggregation not only elevates overall accuracy but also adapts more effectively to cross-dataset heterogeneity than any of the existing baselines.

Table 2: Communication Overhead Results (avg communication time per round in seconds)

Method	Multi-Dataset			Single-Dataset		
	Molecules (7)	Biochem (8)	Mix (13)	IMDB-BINARY (10)	NCII (30)	PROTEINS (10)
FedAvg	0.68	1.07	2.91	0.14	0.49	0.16
GCFL	0.71	1.12	3.03	0.20	0.84	0.21
FedCP	0.60	0.95	2.45	0.12	0.42	0.14
FedProx	0.70	1.09	3.00	0.15	0.50	0.17
FedSage	0.73	1.18	3.20	0.16	0.52	0.18
APPLE	0.75	1.25	3.35	0.17	0.55	0.19
FGSSL	0.85	1.45	4.10	0.19	0.65	0.21
<b>FedSal</b>	1.12	1.78	4.73	0.23	0.86	0.26
FedStar	0.88	1.70	6.70	–	–	–
FedSSP	2.25	4.30	15.50	–	–	–
<b>FedSal+</b>	1.45	2.79	11.69	–	–	–

**Communication Overhead Analysis:** While prior FedGNN studies have not acknowledged communication efficiency in favor of accuracy, we measure average per-round communication time (Table 2) to highlight its importance. Both FedSal and its enhanced variant, FedSal+ incur higher latency than lightweight baselines such as FedAvg, and this gap widens on larger, more heterogeneous datasets. Nevertheless, even at this elevated cost, they still communicate faster than state-of-the-art FedGNN baselines (e.g., FedSSP), preserving a practical edge when bandwidth or response time is critical. Thus, practitioners can attain the accuracy gains of saliency-based aggregation without bearing the highest communication burden among current methods.

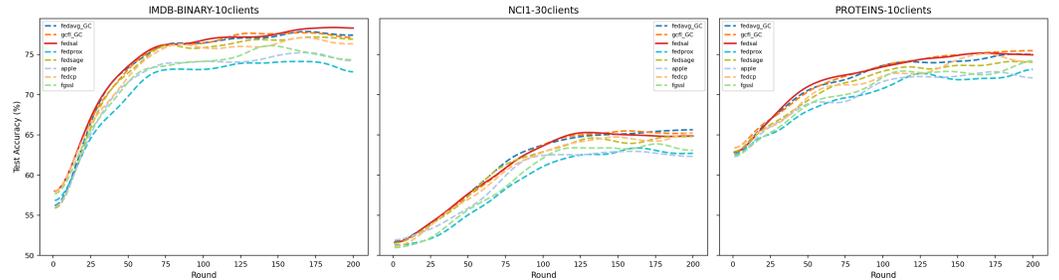


Figure 2: Convergence speed on IMDB-BINARY, NCII-30 and PROTEINS-10 datasets (left to right).

**Analysis of Convergence Speed:** The convergence speed and efficiency of the models are illustrated in Table 3 and Figures 2 and 3. Across single-dataset scenarios, FedSal converges at essentially the same accuracy as mainstream baselines such as FedAvg and GCFL, posting only marginal but consistent gains. In contrast, the advantages of saliency-aware aggregation become clear once cross-dataset heterogeneity is introduced. On the multi-dataset benchmarks, both FedSal and its enhanced variant FedSal+ lead the table, outstripping traditional FL algorithms and other SOTA architectures (e.g., FedStar, FedSSP). The pattern indicates that saliency maps provide marginal benefits when data are IID but yield a tangible edge in non-IID settings by spotlighting informative sub-structures and accelerating effective aggregation. Thus, while the gains in single-dataset scenarios

378  
379  
380  
381  
382  
383  
384  
385  
386  
387

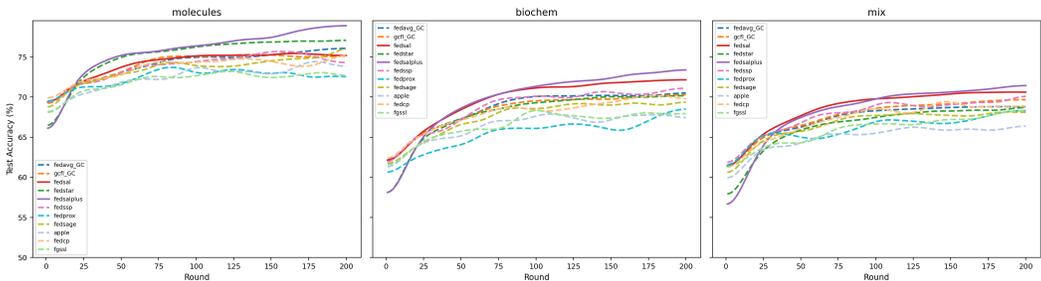


Figure 3: Convergence speed on Molecules, Biochem, Mix datasets (left to right).

Table 3: Accuracy (%) at 50th Communication Round (transposed)

390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402

Method	IMDB-BINARY	NCII	PROTEINS	Molecules	Biochem	Mix
FedAvg	73.34	57.64	70.52	73.34	67.37	66.44
GCFL	73.13	57.39	70.35	73.13	67.19	66.48
FedProx	71.80	55.60	69.00	71.50	66.00	65.80
FedSage	72.60	56.60	70.20	72.10	67.00	66.50
APPLE	72.90	57.50	70.80	72.80	67.80	66.90
FedCP	73.40	57.55	70.95	73.30	68.10	67.00
FGSSL	73.00	57.65	70.85	73.40	68.30	67.10
FedSal	<b>73.55</b>	<b>57.66</b>	<b>71.00</b>	<b>73.55</b>	<b>69.04</b>	<b>67.45</b>
FedStar	–	–	–	73.40	68.20	67.20
FedSSP	–	–	–	73.50	68.55	67.25
FedSal+	–	–	–	<b>73.57</b>	<b>68.63</b>	<b>67.34</b>

403  
404  
405

remain subtle, saliency-based methods maintain, and often extend, their lead as data diversity and complexity grow.

406  
407

### 3.3 HYPERPARAMETER AND ABLATION STUDY

408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421

In Figure 4(a), setting  $\epsilon_{\max}$  too low causes the clustering to ignore meaningful but infrequent gradient updates—preventing the discovery of distinct client behaviors—whereas setting it too high admits spurious, noisy signals that dilute genuine patterns. A mid-range  $\epsilon_{\max}$ , therefore, allows informative deviations to guide cluster formation without letting noise dominate. (b) illustrates the complementary trade-off for  $\epsilon_{\text{mean}}$ : an overly strict mean-threshold treats even small, consistent shifts as convergence—stalling progress, while an overly lenient one smooths away true convergence cues amid stochastic fluctuations. An intermediate value ensures the algorithm remains both responsive to real change and robust to random variation. (c) shows that, in the mixed-domain setting, each successive clustering phase yields a pronounced accuracy gain as primary client subgroups are separated, but returns diminish once the main clusters have formed. Finally, (d) demonstrates that fragmenting clients into too many clusters reduces each group’s sample size below a critical level, undermining the reliability of local model updates. Collectively, these results underscore that both saliency thresholds must be tuned to balance the fundamental noise–signal trade-off in federated clustering.

422

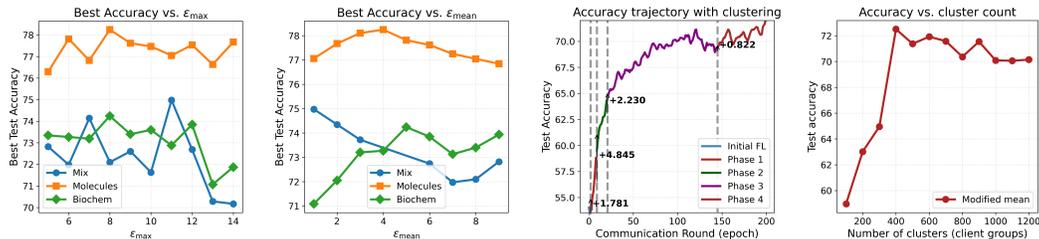
Table 4: Ablation Study of different modules in FedSal compared against accuracy (%).

423  
424  
425  
426  
427  
428  
429

Raw	Sal.	Clust.	IMDB-BINARY	NCII	PROTEINS	Molecules	Biochem	Mix
×	×	✓	76.5±3.2	68.0±1.6	74.0±1.1	74.6±1.1	72.7±0.9	68.9±0.7
×	✓	×	77.8±2.9	69.0±1.6	74.8±1.1	75.4±1.0	70.2±0.8	68.4±0.6
✓	×	✓	75.1±2.7	65.5±1.5	74.9±0.9	75.7±0.9	68.9±0.4	68.9±0.6
×	✓	✓	<b>80.0±2.63</b>	<b>69.9±1.5</b>	<b>75.19±1.0</b>	<b>75.9±1.0</b>	<b>72.8±0.7</b>	<b>70.8±0.5</b>

430  
431

In Table 4, enabling both saliency and clustering yields the strongest performance across all datasets, demonstrating their complementary roles in capturing feature importance and partitioning clients



(a) Effect of  $\epsilon_{\max}$  on test accuracy (b) Effect of  $\epsilon_{\text{mean}}$  on test accuracy (c) Accuracy trajectory across clustering phases (d) Accuracy vs. number of clusters

Figure 4: **Hyper parameter Sensitivity** analysis of FedSal hyper-parameters.

under non-IID conditions. Omitting saliency reduces accuracy notably, despite lower communication overhead, while removing clustering also degrades results, particularly on Biochem and Mix. Using raw gradients matches communication cost but underperforms in accuracy. Together, these results confirm that the combination of saliency maps and dynamic clustering is essential for robust federated learning on heterogeneous graph data.

## 4 DISCUSSION

FedSal and its extension FedSal+ achieve consistently superior accuracy and convergence in both IID and non-IID federated graph learning by leveraging saliency-driven clustering. This saliency-centric approach, however, introduces a measurable communication overhead: transmitting compact saliency summaries alongside model weights increases per-round latency compared to bare-bones methods like FedAvg, yet remains substantially more efficient than spectral or embedding-heavy schemes.

Dynamic clustering underpins nearly half of FedSal’s accuracy gains, but its marginal benefit diminishes once clusters exceed a handful: excessive fragmentation erodes statistical strength and can produce occasional minor instabilities. By tuning adaptive thresholds ( $\epsilon_{\max}$ ,  $\epsilon_{\text{mean}}$ ), FedSal regularizes this saliency noise, preserving smooth, FedAvg-like convergence on IID data while accelerating learning under heterogeneity.

Compared to self-training, which suffers from high variance without aggregation signals, and to gradient-only baselines such as FedCP and FGSSL that either inflate overhead or underperform on homogeneous tasks, the joint mechanism of saliency mapping and dynamic clustering offers the most reliable pathway to robust, cross-domain federated graph learning. Practitioners should balance its modest increase in communication cost against the substantial gains in stability and accuracy when handling heterogeneous graph distributions.

## 5 CONCLUSION

We have presented FedSal, the first federated GNN architecture to exploit saliency activation maps for client clustering and aggregation. By grounding similarity in an interpretable activation space and orchestrating a two-phase clustering–aggregation loop governed by  $\epsilon_{\text{mean}}$  and  $\epsilon_{\max}$ , FedSal mitigates cross-domain conflicts and adapts to evolving heterogeneity without raw-data sharing. Our theoretical results guarantee that saliency distances remain stable under structural or feature perturbations and embed efficiently into Euclidean space. Extensive experiments on thirteen graph classification tasks demonstrate that FedSal delivers modest gains under IID splits and substantial improvements when faced with severe non-IID distributions, outpacing state-of-the-art FedAvg, GCFL, FedStar, and FedSSP. The FedSal+ variant, which enriches node representations with structural embeddings, yields an additional boost in multi-dataset settings, while maintaining a 1 to 3s per-round latency, faster than spectral baselines. Future work will explore Differential privacy, saliency-message compression to reduce bandwidth and validation in large-scale, dynamic hyperparameter tuning. Overall, saliency-guided aggregation, especially when combined with lightweight structural features, offers a practical and high-performance pathway to robust FedGNNs under real-world heterogeneity.

## REFERENCES

- 486  
487  
488 Jinheon Baek, Wonyong Jeong, Jiongdoo Jin, Jaehong Yoon, and Sung Ju Hwang. Personalized  
489 subgraph federated learning. In *International conference on machine learning*, pp. 1396–1415.  
490 PMLR, 2023.
- 491 Jean Bourgain. On lipschitz embedding of finite metric spaces in hilbert space. *Israel Journal of*  
492 *Mathematics*, 52:46–52, 1985.
- 493 Christopher Briggs, Zhong Fan, and Peter Andras. Federated learning with hierarchical clustering of  
494 local updates to improve training on non-iid data. In *2020 international joint conference on neural*  
495 *networks (IJCNN)*, pp. 1–9. IEEE, 2020.
- 497 Chen Chen, Yuchen Liu, Xingjun Ma, and Lingjuan Lyu. Calfat: Calibrated federated adversarial  
498 training with label skewness. *Advances in Neural Information Processing Systems*, 35:3569–3581,  
499 2022a.
- 500 Chen Chen, Lingjuan Lyu, Han Yu, and Gang Chen. Practical attribute reconstruction attack against  
501 federated learning. *IEEE Transactions on Big Data*, 2022b.
- 503 Fei Chen, Mi Luo, Zhenhua Dong, Zhenguo Li, and Xiuqiang He. Federated meta-learning with fast  
504 convergence and efficient communication. *arXiv preprint arXiv:1802.07876*, 2018.
- 506 Fengwen Chen, Guodong Long, Zonghan Wu, Tianyi Zhou, and Jing Jiang. Personalized federated  
507 learning with graph. *arXiv preprint arXiv:2203.00829*, 2022c.
- 508 Mingyang Chen, Wen Zhang, Zonggang Yuan, Yantao Jia, and Huajun Chen. Fede: Embedding  
509 knowledge graphs in federated setting. In *Proceedings of the 10th International Joint Conference*  
510 *on Knowledge Graphs*, pp. 80–88, 2021.
- 511 Yu Chen, Lingfei Wu, and Mohammed J Zaki. Toward subgraph guided knowledge graph question  
512 generation with graph neural networks. *arXiv preprint arXiv:2004.06015*, 2020.
- 514 Zihan Chen, Xingbo Fu, Yuxiao Dong, Jundong Li, and Chunhua Shen. Fedhero: A federated  
515 learning approach for node classification task on heterophilic graphs, 2025. URL [https://](https://arxiv.org/abs/2504.21206)  
516 [arxiv.org/abs/2504.21206](https://arxiv.org/abs/2504.21206).
- 517 Zhiyong Cui, Kristian Henrickson, Ruimin Ke, and Yinhai Wang. Traffic graph convolutional  
518 recurrent neural network: A deep learning framework for network-scale traffic learning and  
519 forecasting. *IEEE Transactions on Intelligent Transportation Systems*, 21(11):4883–4894, 2019.
- 521 Zhehao Dai, Guojiang Shen, Haopeng Yuan, Shangfei Zheng, Yizhou Hu, Jing Du, and Feng Xia.  
522 Towards heterogeneous federated graph learning via structural entropy and prototype aggregation.  
523 *Information Sciences*, 718, 2025. URL [https://www.sciencedirect.com/science/](https://www.sciencedirect.com/science/article/abs/pii/S0020025525004700)  
524 [article/abs/pii/S0020025525004700](https://www.sciencedirect.com/science/article/abs/pii/S0020025525004700).
- 525 Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Personalized federated learning: A meta-  
526 learning approach. *arXiv preprint arXiv:2002.07948*, 2020.
- 527 Hui Fang, Yang Gao, Peng Zhang, Jiangchao Yao, Hongyang Chen, and Haishuai Wang. Large  
528 language models enhanced personalized graph neural architecture search in federated learning.  
529 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 2025. doi: 10.  
530 1609/aaai.v39i16.33814. URL [https://ojs.aaai.org/index.php/AAAI/article/](https://ojs.aaai.org/index.php/AAAI/article/view/33814)  
531 [view/33814](https://ojs.aaai.org/index.php/AAAI/article/view/33814).
- 532 Xingbo Fu, Binchi Zhang, Yushun Dong, Chen Chen, and Jundong Li. Federated graph machine  
533 learning: A survey of concepts, techniques, and applications. *ACM SIGKDD Explorations*  
534 *Newsletter*, 24(2):32–47, 2022.
- 535 Xingbo Fu, Zihan Chen, Binchi Zhang, Chen Chen, and Jundong Li. Federated graph learning with  
536 structure proxy alignment. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge*  
537 *Discovery and Data Mining (KDD '24)*, Barcelona, Spain, 2024. ACM. doi: 10.1145/3637528.  
538 3671717. URL <https://dl.acm.org/doi/10.1145/3637528.3671717>.
- 539

- 540 Xingbo Fu, Zihan Chen, Yinhan He, Song Wang, Binchi Zhang, Chen Chen, and Jundong Li. Virtual  
541 nodes can help: Tackling distribution shifts in federated graph learning. In *Proceedings of the*  
542 *AAAI Conference on Artificial Intelligence*, volume 39, 2025. doi: 10.1609/aaai.v39i16.33830.  
543 URL <https://ojs.aaai.org/index.php/AAAI/article/view/33830>.
- 544 Rufeï Gao, Zhaowei Liu, Chenxi Jiang, Yingjie Wang, Shenqiang Wang, and Pengda Wang. Bi-  
545 fedgnn: Federated graph neural networks framework based on bayesian inference. *Neural Networks*,  
546 169:143–153, 2024. doi: 10.1016/j.neunet.2023.10.024. URL <https://pubmed.ncbi.nlm.nih.gov/37890364/>.
- 547 Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural  
548 message passing for quantum chemistry. In *International conference on machine learning*, pp.  
549 1263–1272. PMLR, 2017.
- 550 Zishan Gu, Ke Zhang, Guangji Bai, Liang Chen, Liang Zhao, and Carl Yang. Dynamic activation of  
551 clients and parameters for federated learning over heterogeneous graphs. ICDE, 2023.
- 552 Jiayan Guo, Shangyang Li, and Yan Zhang. An information theoretic perspective for heterogeneous  
553 subgraph federated learning. In *International Conference on Database Systems for Advanced*  
554 *Applications*, pp. 745–760. Springer, 2023.
- 555 Jenny Hamer, Mehryar Mohri, and Ananda Theertha Suresh. Fedboost: A communication-efficient  
556 algorithm for federated learning. In *International Conference on Machine Learning*, pp. 3973–3983.  
557 PMLR, 2020.
- 558 Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs.  
559 *Advances in neural information processing systems*, 30, 2017.
- 560 Chaoyang He, Keshav Balasubramanian, Emir Ceyani, Carl Yang, Han Xie, Lichao Sun, Lifang He,  
561 Liangwei Yang, S Yu Philip, Yu Rong, et al. Fedgraphnn: A federated learning benchmark system  
562 for graph neural networks. In *ICLR 2021 Workshop on Distributed and Private Machine Learning*  
563 *(DPML)*, 2021.
- 564 Li Huang, Yifeng Yin, Zeng Fu, Shifa Zhang, Hao Deng, and Dianbo Liu. Loadaboost: Loss-based  
565 adaboost federated machine learning with reduced computational complexity on iid and non-iid  
566 intensive care data. *Plos one*, 15(4):e0230706, 2020.
- 567 Wenke Huang, Mang Ye, Zekun Shi, He Li, and Bo Du. Rethinking federated learning with  
568 domain shift: A prototype view. In *2023 IEEE/CVF Conference on Computer Vision and Pattern*  
569 *Recognition (CVPR)*, pp. 16312–16322. IEEE, 2023.
- 570 Wenke Huang, Guancheng Wan, Mang Ye, and Bo Du. Federated graph semantic and structural  
571 learning. *arXiv preprint arXiv:2406.18937*, 2024.
- 572 Eunjeong Jeong, Seungeun Oh, Hyesung Kim, Jihong Park, Mehdi Bennis, and Seong-Lyun Kim.  
573 Communication-efficient on-device machine learning: Federated distillation and augmentation  
574 under non-iid private data. *arXiv preprint arXiv:1811.11479*, 2018.
- 575 Meng Jiang, Taeho Jung, Ryan Karl, and Tong Zhao. Federated dynamic graph neural networks  
576 with secure aggregation for video-based distributed surveillance. *ACM Transactions on Intelligent*  
577 *Systems and Technology (TIST)*, 13(4):1–23, 2022.
- 578 Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin  
579 Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Ad-  
580 vances and open problems in federated learning. *Foundations and trends® in machine learning*,  
581 14(1–2):1–210, 2021.
- 582 Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and  
583 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In  
584 *International conference on machine learning*, pp. 5132–5143. PMLR, 2020.
- 585 Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. Tighter theory for local sgd on identical  
586 and heterogeneous data. In *International Conference on Artificial Intelligence and Statistics*, pp.  
587 4519–4529. PMLR, 2020.

- 594 Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks.  
595 *arXiv preprint arXiv:1609.02907*, 2016.  
596
- 597 Tian Li, Anit Kumar Sahu, Maziar Sanjabi, Manzil Zaheer, Ameet Talwalkar, and Virginia Smith.  
598 On the convergence of federated optimization in heterogeneous networks. *arXiv preprint*  
599 *arXiv:1812.06127*, 2018.
- 600 Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith.  
601 Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems*,  
602 2:429–450, 2020.  
603
- 604 Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated  
605 learning through personalization. In *International conference on machine learning*, pp. 6357–6368.  
606 PMLR, 2021.
- 607 Xianfeng Liang, Shuheng Shen, Jingchang Liu, Zhen Pan, Enhong Chen, and Yifei Cheng. Variance  
608 reduced local sgd with lower communication complexity. *arXiv preprint arXiv:1912.12844*, 2019.  
609
- 610 Jingxin Liu, Jieren Cheng, Renda Han, Wenxuan Tu, Jiaxin Wang, and Xin Peng. Federated  
611 graph-level clustering network. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
612 volume 39, pp. 18870–18878, 2025a. doi: 10.1609/aaai.v39i18.34077. URL [https://ojs.  
613 aaai.org/index.php/AAAI/article/view/34077](https://ojs.aaai.org/index.php/AAAI/article/view/34077).
- 614 Rui Liu, Pengwei Xing, Zichao Deng, Anran Li, Cuntai Guan, and Han Yu. Federated graph neural  
615 networks: Overview, techniques, and challenges. *IEEE Transactions on Neural Networks and*  
616 *Learning Systems*, 2024.
- 617 Yuxuan Liu, Zhiming He, Shuang Wang, Yangyang Wang, Peichao Wang, Zhangshen Huang, and  
618 Qi Sun. Federated subgraph learning via global-knowledge-guided node generation. *Sensors*, 25  
619 (7):2240, 2025b. doi: 10.3390/s25072240. URL [https://www.mdpi.com/1424-8220/  
620 25/7/2240](https://www.mdpi.com/1424-8220/25/7/2240).
- 621 Guannan Lou, Yuze Liu, Tiehua Zhang, and Xi Zheng. Stfl: A temporal-spatial federated learning  
622 framework for graph neural networks. *arXiv preprint arXiv:2111.06750*, 2021.  
623
- 624 Jun Luo and Shandong Wu. Adapt to adaptation: Learning personalization for cross-silo federated  
625 learning. In *IJCAI: proceedings of the conference*, volume 2022, pp. 2166, 2022.  
626
- 627 Mi Luo, Fei Chen, Dapeng Hu, Yifan Zhang, Jian Liang, and Jiashi Feng. No fear of heterogeneity:  
628 Classifier calibration for federated learning with non-iid data. *Advances in Neural Information*  
629 *Processing Systems*, 34:5972–5984, 2021.
- 630 Lingjuan Lyu, Han Yu, Xingjun Ma, Chen Chen, Lichao Sun, Jun Zhao, Qiang Yang, and S Yu Philip.  
631 Privacy and robustness in federated learning: Attacks and defenses. *IEEE transactions on neural*  
632 *networks and learning systems*, 2022.  
633
- 634 Chengyuan Mai, Tianchi Liao, Chuan Chen, and Zibin Zheng. Fgtl: Federated graph transfer  
635 learning for node classification. *ACM Transactions on Knowledge Discovery from Data*, 2024. doi:  
636 10.1145/3699962. URL <https://dl.acm.org/doi/10.1145/3699962>.
- 637 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.  
638 Communication-efficient learning of deep networks from decentralized data. In *Artificial intelli-*  
639 *gence and statistics*, pp. 1273–1282. PMLR, 2017.  
640
- 641 H McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and B Aguera y Arcas. Communication-  
642 efficient learning of deep networks from decentralized data. arxiv e-prints. *arXiv preprint*  
643 *arXiv:1602.05629*, 2016.
- 644 Kurt Mehlhorn and Christian Uhrig. The minimum cut algorithm of stoer and wagner. 1995.  
645
- 646 Chuizheng Meng, Sirisha Rambhatla, and Yan Liu. Cross-node federated graph neural network  
647 for spatio-temporal data modeling. In *Proceedings of the 27th ACM SIGKDD conference on*  
*knowledge discovery & data mining*, pp. 1202–1211, 2021.

- 648 Wenjie Pei, Weina Xu, Zongze Wu, Weichao Li, Jinfan Wang, Guangming Lu, and Xiangrong Wang.  
649 Saliency-aware regularized graph neural network, 2024. URL [https://arxiv.org/abs/  
650 2401.00755](https://arxiv.org/abs/2401.00755).
- 651 Felix Sattler, Klaus-Robert Müller, and Wojciech Samek. Clustered federated learning: Model-  
652 agnostic distributed multitask optimization under privacy constraints. *IEEE transactions on neural  
653 networks and learning systems*, 32(8):3710–3722, 2020.
- 654  
655 Abdullah Abdul Sattar Shaikh and Saeed Samet. Federated graph neural networks in non-iid scenar-  
656 ios—a comprehensive survey. *Neurocomputing*, 2025. URL [https://www.sciencedirect.  
657 com/science/article/abs/pii/S0925231225016790](https://www.sciencedirect.com/science/article/abs/pii/S0925231225016790). Article 131007.
- 658  
659 Toyotaro Suzumura, Yi Zhou, Natahalie Baracaldo, Guangnan Ye, Keith Houck, Ryo Kawahara, Ali  
660 Anwar, Lucia Larise Stavarache, Yuji Watanabe, Pablo Loyola, et al. Towards federated graph  
661 learning for collaborative financial crimes detection. *arXiv preprint arXiv:1909.12946*, 2019.
- 662  
663 Canh T Dinh, Nguyen Tran, and Josh Nguyen. Personalized federated learning with moreau envelopes.  
664 *Advances in neural information processing systems*, 33:21394–21405, 2020.
- 665  
666 Yue Tan, Guodong Long, Lu Liu, Tianyi Zhou, Qinghua Lu, Jing Jiang, and Chengqi Zhang. Fedproto:  
667 Federated prototype learning across heterogeneous clients. In *Proceedings of the AAAI Conference  
668 on Artificial Intelligence*, volume 36, pp. 8432–8440, 2022a.
- 669  
670 Yue Tan, Guodong Long, Jie Ma, Lu Liu, Tianyi Zhou, and Jing Jiang. Federated learning from  
671 pre-trained models: A contrastive learning approach. *Advances in neural information processing  
672 systems*, 35:19332–19344, 2022b.
- 673  
674 Yue Tan, Yixin Liu, Guodong Long, Jing Jiang, Qinghua Lu, and Chengqi Zhang. Federated learning  
675 on non-iid graphs via structural knowledge sharing. In *Proceedings of the AAAI conference on  
676 artificial intelligence*, volume 37, pp. 9953–9961, 2023.
- 677  
678 Zihan Tan, Guancheng Wan, Wenke Huang, and Mang Ye. Fedssp: Federated graph learning with  
679 spectral knowledge and personalized preference. *arXiv preprint arXiv:2410.20105*, 2024.
- 680  
681 Zihan Tan, Guancheng Wan, Wenke Huang, He Li, Guibin Zhang, Carl Yang, and Mang  
682 Ye. Fedspa: Generalizable federated graph learning under homophily heterogeneity. In  
683 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,  
684 2025. URL [https://openaccess.thecvf.com/content/CVPR2025/html/Tan\\_  
685 FedSPA\\_Generalizable\\_Federated\\_Graph\\_Learning\\_under\\_Homophily\\_  
686 Heterogeneity\\_CVPR\\_2025\\_paper.html](https://openaccess.thecvf.com/content/CVPR2025/html/Tan_FedSPA_Generalizable_Federated_Graph_Learning_under_Homophily_Heterogeneity_CVPR_2025_paper.html).
- 687  
688 Binghui Wang, Ang Li, Meng Pang, Hai Li, and Yiran Chen. Graphfl: A federated learning framework  
689 for semi-supervised node classification on graphs. In *2022 IEEE International Conference on Data  
690 Mining (ICDM)*, pp. 498–507. IEEE, 2022a.
- 691  
692 Hengyi Wang, Weiyang Xie, Jitao Ma, Daixun Li, and Yunsong Li. FedSLS: Exploring feder-  
693 ated aggregation in saliency latent space. In *ACM Multimedia 2024*, 2024. URL [https:  
694 //openreview.net/forum?id=epLGidFcqi](https://openreview.net/forum?id=epLGidFcqi).
- 695  
696 Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. Tackling the objective  
697 inconsistency problem in heterogeneous federated optimization. *Advances in neural information  
698 processing systems*, 33:7611–7623, 2020.
- 699  
700 Shiqiang Wang, Tiffany Tuor, Theodoros Salonidis, Kin K Leung, Christian Makaya, Ting He, and  
701 Kevin Chan. Adaptive federated learning in resource constrained edge computing systems. *IEEE  
702 journal on selected areas in communications*, 37(6):1205–1221, 2019.
- 703  
704 Zhen Wang, Weirui Kuang, Yuexiang Xie, Liuyi Yao, Yaliang Li, Bolin Ding, and Jingren Zhou.  
705 Federatedscope-gnn: Towards a unified, comprehensive and efficient package for federated graph  
706 learning. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and  
707 Data Mining (KDD '22)*, Washington, DC, USA, 2022b. ACM. doi: 10.1145/3534678.3539112.  
708 URL <https://dl.acm.org/doi/10.1145/3534678.3539112>.

- 702 Chuhan Wu, Fangzhao Wu, Lingjuan Lyu, Tao Qi, Yongfeng Huang, and Xing Xie. A federated  
703 graph neural network framework for privacy-preserving personalization. *Nature Communications*,  
704 13(3091), 2022. doi: 10.1038/s41467-022-30714-9. URL [https://www.nature.com/  
705 articles/s41467-022-30714-9](https://www.nature.com/articles/s41467-022-30714-9).
- 706 Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. A  
707 comprehensive survey on graph neural networks. *IEEE transactions on neural networks and  
708 learning systems*, 32(1):4–24, 2020.
- 709 Han Xie, Jing Ma, Li Xiong, and Carl Yang. Federated graph classification over non-iid graphs.  
710 *Advances in neural information processing systems*, 34:18839–18852, 2021.
- 711 Han Xie, Li Xiong, and Carl Yang. Federated node classification over graphs with latent link-type  
712 heterogeneity. In *Proceedings of the ACM Web Conference 2023*, pp. 556–566, 2023.
- 713 Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie  
714 Jegelka. Representation learning on graphs with jumping knowledge networks. In *International  
715 conference on machine learning*, pp. 5453–5462. PMLR, 2018.
- 716 Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. Federated machine learning: Concept and  
717 applications. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 10(2):1–19, 2019.
- 718 Dezhong Yao, Wanning Pan, Yutong Dai, Yao Wan, Xiaofeng Ding, Chen Yu, Hai Jin, Zheng Xu, and  
719 Lichao Sun. Fedgkd: Toward heterogeneous federated learning via global knowledge distillation.  
720 *IEEE Transactions on Computers*, 73(1):3–17, 2024. doi: 10.1109/TC.2023.3315066.
- 721 Hao Yu, Sen Yang, and Shenghuo Zhu. Parallel restarted sgd with faster convergence and less  
722 communication: Demystifying why model averaging works for deep learning. In *Proceedings of  
723 the AAAI conference on artificial intelligence*, volume 33, pp. 5693–5700, 2019.
- 724 Wentao Yu, Shuo Chen, Yongxin Tong, Tianlong Gu, and Chen Gong. Modeling inter-intra  
725 heterogeneity for graph federated learning. In *Proceedings of the AAAI Conference on Artificial  
726 Intelligence*, volume 39, pp. 22236–22244, 2025. doi: 10.1609/aaai.v39i21.34378. URL  
727 <https://ojs.aaai.org/index.php/AAAI/article/view/34378>.
- 728 Xiaotong Yuan and Ping Li. On convergence of fedprox: Local dissimilarity invariant bounds, non-  
729 smoothness and beyond. *Advances in Neural Information Processing Systems*, 35:10752–10765,  
730 2022.
- 731 Huanding Zhang, Tao Shen, Fei Wu, Mingyang Yin, Hongxia Yang, and Chao Wu. Federated graph  
732 learning—a position paper. *arXiv preprint arXiv:2105.11099*, 2021a.
- 733 Jianqing Zhang, Yang Hua, Hao Wang, Tao Song, Zhengui Xue, Ruhui Ma, and Haibing Guan.  
734 Fedcp: Separating feature information for personalized federated learning via conditional policy.  
735 In *Proceedings of the 29th ACM SIGKDD conference on knowledge discovery and data mining*,  
736 pp. 3249–3261, 2023.
- 737 Ke Zhang, Carl Yang, Xiaoxiao Li, Lichao Sun, and Siu Ming Yiu. Subgraph federated learning with  
738 missing neighbor generation. *Advances in Neural Information Processing Systems*, 34:6671–6682,  
739 2021b.
- 740 Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated  
741 learning with non-iid data. *arXiv preprint arXiv:1806.00582*, 2018.
- 742 Longfei Zheng, Jun Zhou, Chaochao Chen, Bingzhe Wu, Li Wang, and Benyu Zhang. Afsfgnn:  
743 Automated separated-federated graph neural network. *Peer-to-Peer Networking and Applications*,  
744 14(3):1692–1704, 2021.
- 745 Wei Zhu, Jiebo Luo, and Andrew D White. Federated learning of molecular properties with graph  
746 neural networks in a heterogeneous setting. *Patterns*, 3(6), 2022.
- 747  
748  
749  
750  
751  
752  
753  
754  
755

## A APPENDIX / SUPPLEMENTAL MATERIAL

Due to the page limitations, this section provides further theoretical and experimental details for the FedSal framework.

### A.1 CLUSTERING MECHANISM

The top-down clustering strategy referenced in Section 2.1.2 operates as follows. Let a cluster  $C \subseteq \{1, \dots, N\}$  contain  $|C|$  clients, each transmitting a saliency map update  $\Delta S_i \in \mathbb{R}^d$ . Define

$$\begin{aligned}\epsilon_{\text{mean}}(C) &= \frac{1}{|C|} \sum_{i \in C} \|\Delta S_i\|, \\ \epsilon_{\text{max}}(C) &= \max_{i \in C} \|\Delta S_i\|.\end{aligned}$$

When

$$\epsilon_{\text{mean}}(C) < \epsilon_{\text{mean}} \quad \text{and} \quad \epsilon_{\text{max}}(C) > \epsilon_{\text{max}},$$

we conclude that the cluster converges on average (small mean update), yet contains at least one outlier client (large max update). Formally,

$$\epsilon_{\text{mean}}(C) < \epsilon_{\text{mean}} \quad \text{and} \quad \epsilon_{\text{max}}(C) > \epsilon_{\text{max}} \implies (\text{split } C).$$

In this event, we *partition*  $C$  into two sub-clusters, denoted  $C_1$  and  $C_2$ .

To find the partition, we construct a fully connected graph  $G = (V, E)$  where  $V = C$  and  $E = \{(i, j) : i, j \in C\}$ . Each edge  $(i, j)$  is assigned a weight

$$\alpha_{ij} = \frac{\Delta S_i \cdot \Delta S_j}{\|\Delta S_i\| \|\Delta S_j\|},$$

which represents the directional alignment of the saliency-map updates. Large  $\alpha_{ij}$  indicates high alignment in salient feature gradients, suggesting that clients  $i$  and  $j$  have similar data distributions or tasks.

Since a minimum-cut algorithm would naively remove edges of largest total weight, we convert similarities into dissimilarities:

$$w_{ij} = 1 - \alpha_{ij}.$$

A lower  $w_{ij}$  therefore implies greater similarity. We then apply the Stoer–Wagner minimum-cut algorithm Mehlhorn & Urig (1995) to minimize

$$\sum_{i \in C_1} \sum_{j \in C_2} w_{ij} = \sum_{i \in C_1} \sum_{j \in C_2} (1 - \alpha_{ij}),$$

thereby favoring that highly similar clients ( $\alpha_{ij} \approx 1$ ) remain together. The result is a bi-partition  $\{C_1, C_2\} \subseteq C$  that groups mutually similar clients.

After partitioning, we recursively test each new sub-cluster  $C_k \in \{C_1, C_2\}$  using the same thresholds. If

$$\epsilon_{\text{mean}}(C_k) < \epsilon_{\text{mean}} \quad \text{and} \quad \epsilon_{\text{max}}(C_k) > \epsilon_{\text{max}},$$

we again split  $C_k$ , repeating until no cluster satisfies the splitting condition.

During subsequent communication rounds, each cluster  $C_k$  updates its model using only the aggregated updates from its member clients. Formally, let  $\theta_{C_k}$  denote the parameters for cluster  $C_k$ . Each client  $i \in C_k$  solves

$$\theta_i^{(t)} = \arg \min_{\theta} \mathcal{L}(\theta; D_i)$$

and transmits its update  $\Delta \theta_i^{(t)}$  (and corresponding  $\Delta S_i$ ) to the cluster aggregator, which performs

$$\theta_{C_k}^{(t+1)} = \theta_{C_k}^{(t)} + \frac{1}{|C_k|} \sum_{i \in C_k} \Delta \theta_i^{(t)}.$$

This localized aggregation creates a specialized model for each cluster’s data characteristics, mitigating performance degradation from heterogeneous distributions. Once no further splits occur, the final clusters  $\{C_1, \dots, C_K\}$  stabilize and each sub-model  $\theta_{C_k}$  converges on its subset of clients.

In practice, the thresholds  $\epsilon_{\text{mean}}$  and  $\epsilon_{\text{max}}$  are tuned to reflect acceptable intra-cluster variability. For highly non-IID data, clustering yields markedly better performance than a single global model, at the cost of extra Stoer–Wagner computations—often justified by faster convergence and reduced variance across heterogeneous clients.

## A.2 PRELIMINARIES

### A.2.1 GRAPH NEURAL NETWORKS

GNNs represent an advanced paradigm in artificial neural networks, specifically designed for processing data in network or graph structures. Let  $G = (V, E)$  be a graph where  $V$  is a set of nodes,  $E$  is a set of edges, and  $X$  denotes node features. GNNs aim to learn node-level representations  $h_v$  for  $v \in V$  and/or graph-level representations  $h_G$  for the entire graph  $G$ . While GNNs have various types Kipf & Welling (2016); Hamilton et al. (2017); Gilmer et al. (2017); Xu et al. (2018) The fundamental operation in GNNs is based on message passing and neighbourhood aggregation. The process consists of the following stages:

**Message Passing:** In this stage, each node  $v$  sends messages to its neighbours. The messages contain information about the node’s features, and the process can be mathematically described as:

$$m_{uv}^{(l)} = \text{MSG}^{(l)}(h_u^{(l)}, h_v^{(l)}, e_{uv}), \quad (1)$$

where  $m_{uv}^{(l)}$  is the message sent from node  $u$  to node  $v$  at layer  $l$ ,  $h_u^{(l)}$  and  $h_v^{(l)}$  are the representations of nodes  $u$  and  $v$  at layer  $l$ , and  $e_{uv}$  represents the edge features between nodes  $u$  and  $v$ .

**Neighborhood Aggregation:** After receiving messages from its neighbours, each node aggregates these messages to update its representation. This can be formalized as:

$$a_v^{(l)} = \text{AGGREGATE}^{(l)}(\{m_{uv}^{(l)} : u \in N(v)\}), \quad (2)$$

where  $a_v^{(l)}$  is the aggregated message for node  $v$  at layer  $l$ , and  $N(v)$  is the set of neighbors of node  $v$ .

**Update:** The node updates its representation based on the aggregated message:

$$h_v^{(l+1)} = \text{UPDATE}^{(l)}(h_v^{(l)}, a_v^{(l)}), \quad (3)$$

where  $h_v^{(l+1)}$  is the updated representation of node  $v$  at layer  $l + 1$ . These steps are repeated for each layer in the GNN, allowing nodes to progressively incorporate information from further distances in the graph. Graph-level representation  $h_G$  can be obtained by aggregating node representations using readout functions like mean pooling or sum pooling:

$$h_G = \text{READOUT}(\{h_v : v \in V\}). \quad (4)$$

### A.2.2 FEDERATED LEARNING

FL is a distributed machine learning approach that enables multiple entities to collaboratively train a model without sharing raw data. Consider  $M$  clients, each with a private dataset  $D_m$ . The global objective is to minimize:

$$\min_{\theta_1, \theta_2, \dots, \theta_M} \frac{1}{M} \sum_{m=1}^M \frac{|D_m|}{N} L_m(\theta_m; D_m), \quad (5)$$

where  $N$  is the total number of instances, and  $L_m$  and  $\theta_m$  are the loss function and model parameters of client  $m$ . The FL process involves several key steps:

**Initialization:** A global model  $\theta_0$  is initialized by the central server without pre-existing knowledge.

**Local Training:** Each client  $m$  downloads the current global model  $\theta_t$  and trains it on its local dataset  $D_m$  for a specified number of local epochs  $E$ . The local training updates the model parameters using gradient descent:

$$\theta_m^{(t+1)} = \theta_t - \eta \nabla L_m(\theta_t; D_m), \quad (6)$$

864 where  $\eta$  is the learning rate.

865 **Aggregation:** The central server aggregates the received model parameters from all clients to update  
866 the global model. The most commonly used aggregation algorithm in FL is FedAvg, as described by  
867 McMahan et al. (2017) (Algorithm 1 in Appendix Section A.6). The FedAvg algorithm computes a  
868 simple average of the model parameters from all participating clients:  
869

$$870 \theta_{t+1} = \frac{1}{M} \sum_{m=1}^M \theta_m^{(t+1)}, \quad (7)$$

872 where  $M$  is the number of clients that have contributed to the aggregation process. This method  
873 ensures that each client, regardless of the size of their dataset, contributes equally to the global model.  
874 The aggregation step effectively combines the local updates from different clients, integrating the  
875 knowledge learned from their individual datasets.

876 **Communication Rounds:** This process of local training and aggregation is repeated over multiple  
877 communication rounds. At each round  $t$ , the global model  $\theta_t$  is progressively refined, improving its  
878 performance on the overall task. Each round consists of: 1. The server sending the current global  
879 model to the clients. 2. Clients performing local training and sending the updated model parameters  
880 back to the server. 3. The server aggregating the updates to form a new global model. By iterating  
881 through these communication rounds, FL allows the global model to learn from diverse data sources  
882 without conflict of knowledge among them.

### 884 A.2.3 SALIENCY MAPS

885 Saliency maps are a technique for interpreting neural networks by highlighting the importance of  
886 input features. Given an input  $\mathbf{x}$  and a neural network model  $f(\mathbf{x}; \theta)$  with parameters  $\theta$ , the saliency  
887  $S(\mathbf{x})$  of each input feature  $x_i$  is quantified by the gradient of the loss function  $\mathcal{L}$  with respect to  $x_i$ .  
888 Formally, the saliency map  $S$  for an input  $\mathbf{x}$  is defined as:

$$889 S(\mathbf{x}) = \left| \frac{\partial \mathcal{L}(f(\mathbf{x}; \theta), y)}{\partial \mathbf{x}} \right|, \quad (8)$$

892 where  $\mathcal{L}$  is the loss function,  $f(\mathbf{x}; \theta)$  is the model’s prediction, and  $y$  is the true label. The gradient  
893  $\frac{\partial \mathcal{L}(f(\mathbf{x}; \theta), y)}{\partial \mathbf{x}}$  captures the sensitivity of the loss with respect to the input features, indicating the  
894 importance of each feature in the model’s decision. The absolute value is taken to capture the  
895 magnitude of importance, regardless of the direction of the gradient. Higher saliency values indicate  
896 features that are more influential in the model’s output. Further Discussion on Saliency Maps can be  
897 found in Appendix Section A.3.

### 898 A.3 ADDITIONAL DISCUSSION ON SALIENCY MAPS

900 Saliency maps are widely used in neural networks to visualize and interpret the importance of input  
901 features. The concept of saliency can be understood with a simple example: consider an image  
902 classification task where a neural network is trained to identify cats and dogs in images. For the input  
903 image  $\mathbf{x}$  and model parameters  $\theta$ , the saliency  $S(\mathbf{x})$  of each pixel  $x_i$  is determined by the gradient  
904 of the loss function  $\mathcal{L}$  with respect to  $x_i$ . Formally, we define the saliency map  $S$  for an input  $\mathbf{x}$  as  
905 follows:

$$906 S(\mathbf{x}) = \left| \frac{\partial \mathcal{L}(f(\mathbf{x}; \theta), y)}{\partial \mathbf{x}} \right|, \quad (9)$$

910 where  $\mathcal{L}$  represents the loss function,  $f(\mathbf{x}; \theta)$  is the model’s prediction, and  $y$  is the true label of the  
911 image. This gradient  $\frac{\partial \mathcal{L}(f(\mathbf{x}; \theta), y)}{\partial \mathbf{x}}$  captures how sensitive the loss is to changes in the input features,  
912 thus highlighting the significance of each feature in the decision-making process of the model. Higher  
913 saliency values signify features that are more influential in affecting the model’s output.

914 Saliency maps, such as the one shown in Figure 5 on the right, visually represent the regions within  
915 an image that most strongly influence the model’s output. This technique is useful in understanding  
916 which features the model deems significant, thus providing insights into the model’s decision-making  
917 process. For instance, in this example, the saliency map reveals that the model focuses predominantly  
on and around the face and fur cat, areas typically critical for animal recognition tasks.

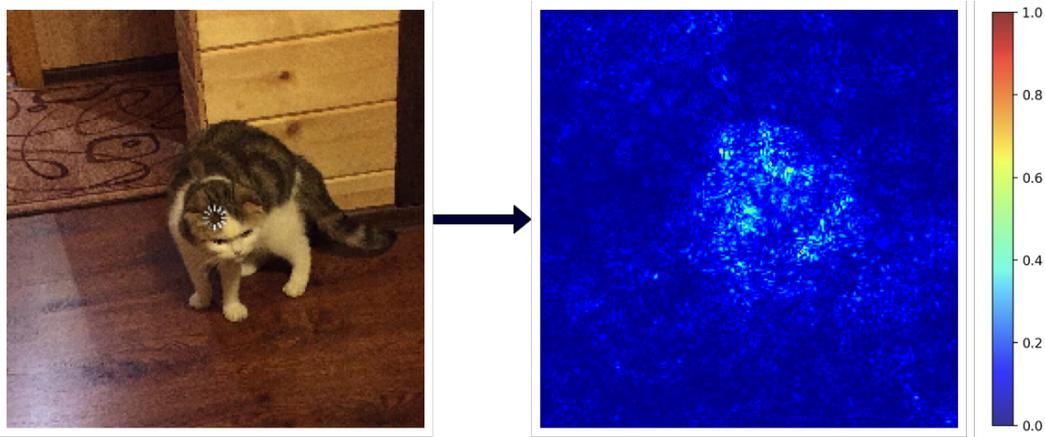


Figure 5: Saliency example

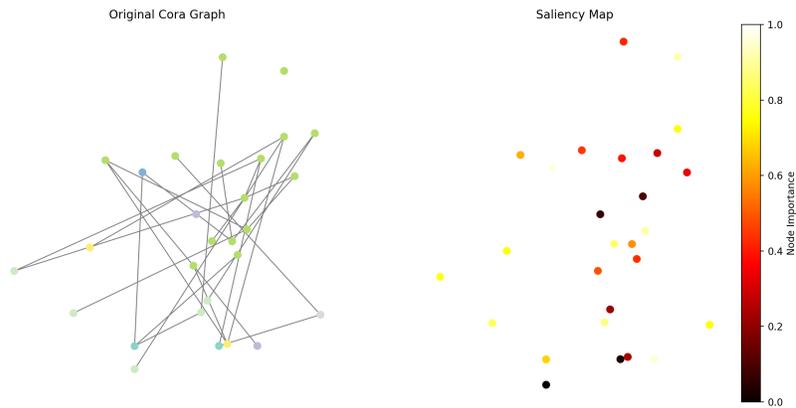


Figure 6: Graph saliency

### A.3.1 SALIENCY IN GRAPH NEURAL NETWORKS

In the context of GNNs, saliency maps can be used to identify the most influential nodes, edges, or features within a graph that contribute to the model’s predictions. Consider a social network graph where nodes represent users and edges represent relationships between users. for a graph  $G = (V, E, X)$ , where  $V$  are the nodes,  $E$  the edges, and  $X$  the node features, the saliency map  $S$  for a GNN model  $f(G; \theta)$  can be similarly defined:

$$S(G) = \left| \frac{\partial \mathcal{L}(f(G; \theta), y)}{\partial X} \right|, \quad (10)$$

Here,  $\frac{\partial \mathcal{L}(f(G; \theta), y)}{\partial X}$  denotes the gradient of the model’s loss function with respect to the node features  $X$ , offering insights into which features are critical in determining the model’s predictions. The gradient provides insights into the importance of different node features in determining the model’s predictions. The provided Figure 6 shows the concept of saliency maps in GNNs using the Cora dataset. On the left, we see the original Cora graph, where nodes and edges are depicted without any indication of their relative importance. On the right, the saliency map highlights the nodes’ importance, which is crucial for the model’s predictions.

In the saliency map, nodes are colour-coded to indicate their significance, with the colour intensity corresponding to the magnitude of their importance. Darker colours (red and black) represent higher importance, whereas lighter colours (yellow) indicate lower importance. This visualization helps in identifying which nodes and their features contribute most to the model’s output. For instance, nodes with higher importance (darker colours) may represent critical information hubs or influential entities within the graph, significantly impacting the GNN’s classification decisions.

By analyzing saliency maps, researchers and practitioners can gain valuable insights into the inner workings of GNNs, understanding how the model interprets the structure and features of graph data. This can aid in model debugging, feature engineering, and improving the interpretability of GNN-based models. Additionally, saliency maps can be used to refine the model by focusing on the most influential nodes and edges, leading to better performance and more robust predictions.

#### A.4 TECHNICAL DETAILS ON FEDSAL+

Influenced by the FedStar framework Tan et al. (2023), we experiment with **FedSal+**, integrating advanced encoding strategies, specifically positional and random walk encodings, to enhance the structural detail in FL clients. Most details, hyperparameters, and model architecture were taken from this reference to guarantee performance. These enhancements aim to refine client-side feature representations, thereby improving the efficacy of FedSal’s saliency-based model aggregation. By deepening the feature set available for each node, FedSal+ seeks to enrich the saliency maps produced, refining the focus on pivotal features during FL. This integration tests the modularity and heightens the compatibility of FedSal with orthogonal methods, such as structural embeddings. The key components of FedSal+ include:

- **Degree-Based Structure Embedding (DSE):** Utilizes vertex degrees in a one-hot encoding format to encapsulate local structural knowledge. This method, while simple, effectively captures fundamental geometric properties of nodes and ensures computational efficiency.
- **Random Walk-Based Structure Embedding (RWSE):** Leverages the random walk diffusion process to assess global structural patterns. This embedding evaluates the probability distribution of node connectivity over several steps, enriching the node’s structural context.

Combining these embeddings results in a comprehensive structural descriptor for each node:

$$s_v = \text{concat}[s_{\text{DSE},v}, s_{\text{RWSE},v}]$$

This concatenated embedding  $s_v$  provides a holistic view of both local and global structural knowledge, significantly enriching the node’s feature set.

The primary aim of FedSal+ is to determine whether this enriched feature set can work with FedSal’s saliency map-driven model aggregation to bolster the framework’s effectiveness, particularly in handling non-IID data across federated settings. We hypothesize that the enriched, refined feature representations facilitated by FedSal+ will lead to more precise and informative saliency maps, thereby improving overall model performance in federated environments. This hypothesis sets the stage for our experimental investigations, aiming to validate the improved performance metrics that FedSal+ could potentially deliver by effectively utilizing advanced structural embeddings.

#### A.5 EXPERIMENTAL DETAILS

Table 5 shows the feature details of the dataset utilized.

##### A.5.1 DATASET DETAILS

##### A.5.2 HYPERPARAMETERS TESTED

We carried out a focused grid search on critical hyper-parameters on the validation dataset, while less sensitive parameters were fixed. The grid search covered:

- Learning rate:  $\{0.005, 0.001, 0.0005, 0.0001\}$
- Weight decay:  $\{0.0007, 0.0005, 0.0003, 0.0001\}$

Table 5: Dataset Details

Dataset Name	No # Graphs	Avg. # Nodes	Avg. #Edges	# Labels	No of Node Features
MUTAG	188	17.93	19.79	2	7
BZR	405	35.75	38.36	2	53
COX2	467	41.22	43.45	2	35
DHFR	467	42.43	44.54	2	53
PTC_MR	344	14.29	14.69	2	18
AIDS	2000	15.69	16.20	2	38
NCII	4110	29.87	32.30	2	37
ENZYMES	600	32.63	62.14	6	3
DD	1178	284.32	715.66	2	89
PROTEINS	1113	39.06	72.82	2	3
COLLAB	5000	74.49	2457.78	3	1 (degree)
IMDB-BINARY	1000	19.77	96.53	2	1 (degree)
IMDB-MULTI	1500	13.00	65.94	3	1 (degree)

- $\epsilon_{\text{mean}}: \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$
- $\epsilon_{\text{max}}: \{5, 6, 7, 8, 9, 10, 11, 12, 13, 14\}$

### A.5.3 MODEL ARCHITECTURE

Table 6: Layer-by-layer configuration of the GIN model used for graph-based learning tasks.

Layer	Type	Details
1	FC	$n_{\text{feat}} \rightarrow 64$
2-4	GIN + ReLU + Dropout	$128 \rightarrow 64$ , dropout = 0.5
5	Pooling	Global Additive Pooling
6	FC	$128 \rightarrow 64$
7	FC + ReLU + Dropout	$64 \rightarrow 64$ , dropout = 0.5
8	FC	$64 \rightarrow n_{\text{class}}$

Table 7: Layer-by-layer configuration of the FedSal+ model with.

Layer	Type	Details
1	FC (Input)	$n_{\text{feat}} = 64 \rightarrow 128$
2	FC (Embedding of $n_{\text{se}}$ )	$n_{\text{se}} = 32 \rightarrow 128$
3-5	GIN + ReLU + Dropout	$256 \rightarrow 128$ , dropout = 0.5
3-5	GCN (for $n_{\text{se}}$ features)	$128 \rightarrow 128$
3-5	Concatenation (x, s)	Concatenation of feature $x$ and $s$
6	Pooling	Global Additive Pooling
7	FC + ReLU	$128 \rightarrow 128$
8	FC + Dropout	$128 \rightarrow 2$ , dropout = 0.5

The GIN model used in FedSal and FedSal+ employ several layers to optimize the processing of graph data. The model configuration is detailed in Table 6 and Table 7 respectively.

## A.6 ALGORITHMS

---

**Algorithm 1** FedSal: Federated Learning with Saliency Aggregation for Graph Neural Networks

---

**Require:** Initial global model  $\theta$   
**Require:** Local datasets  $\{D_i\}_{i=1}^N$   
**Require:** Number of communication rounds  $T$   
**Require:** Clustering thresholds  $\epsilon_{\text{mean}}, \epsilon_{\text{max}}$

- 1: Initialize clusters  $\mathcal{C} \leftarrow \{C_1\}$
- 2: Initialize  $S_{\text{prev}}[i] \leftarrow 0$  for all  $i$
- 3: **for** each round  $t = 1, 2, \dots, T$  **do**
- 4:   Broadcast global model  $\theta$  to all clients each client  $i = 1, \dots, N$  **in parallel**
- 5:    $(\theta_i, S_i) \leftarrow \text{LOCALUPDATE}(i, \theta)$
- 6:    $\Delta S_i \leftarrow S_i - S_{\text{prev}}[i]$
- 7:    $S_{\text{prev}}[i] \leftarrow S_i$
- 8:   **if**  $t > 20$  **then**
- 9:      $\mathcal{C} \leftarrow \text{CLUSTERCLIENTS}(\{\Delta S_i\}, \mathcal{C}, \epsilon_{\text{mean}}, \epsilon_{\text{max}})$
- 10:   **end if**
- 11:    $\theta \leftarrow \text{AGGREGATEMODELS}(\mathcal{C}, \{\theta_i\})$
- 12: **end for**

---



---

**Algorithm 2** Dynamic Client Clustering

---

- 1: **function** CLUSTERCLIENTS( $\{\Delta S_i\}_{i=1}^N, \mathcal{C}, \epsilon_{\text{mean}}, \epsilon_{\text{max}}$ )
- 2:   **for** each cluster  $C_k \in \mathcal{C}$  **do**
- 3:      $\delta_{\text{mean}}^k \leftarrow \frac{1}{|C_k|} \sum_{i \in C_k} \|\Delta S_i\|$
- 4:      $\delta_{\text{max}}^k \leftarrow \max_{i \in C_k} \|\Delta S_i\|$
- 5:     **if**  $\delta_{\text{mean}}^k < \epsilon_{\text{mean}}$  **and**  $\delta_{\text{max}}^k > \epsilon_{\text{max}}$  **then**
- 6:       Compute  $\alpha_{ij} \leftarrow \frac{\Delta S_i \cdot \Delta S_j}{\|\Delta S_i\| \|\Delta S_j\|}$  for all  $i, j \in C_k$
- 7:        $w_{ij} \leftarrow 1 - \alpha_{ij}$
- 8:       Apply Stoer–Wagner min-cut on graph  $(C_k, \{w_{ij}\})$  to obtain  $C_{k1}, C_{k2}$
- 9:       Replace  $C_k$  in  $\mathcal{C}$  with  $C_{k1}, C_{k2}$
- 10:    **end if**
- 11:   **end for**
- 12:   **return**  $\mathcal{C}$
- 13: **end function**

---

**Algorithm 3** Local Update and Saliency Aggregation

---

**Require:** Client index  $i$ , global model  $\theta$   
**Require:** Local dataset  $D_i$ , local epochs  $E$ , learning rate  $\eta$

```

1: function LOCALUPDATE( $i, \theta$ )
2:    $\theta_i \leftarrow \theta$ 
3:   for epoch  $e = 1, \dots, E$  do
4:      $\theta_i \leftarrow \theta_i - \eta \nabla_{\theta_i} \mathcal{L}(\theta_i; D_i)$ 
5:   end for
6:    $S_i \leftarrow 0$ 
7:   for each  $(\mathbf{x}, y) \in D_i$  do
8:      $S(\mathbf{x}) \leftarrow \left| \frac{\partial \mathcal{L}(f(\mathbf{x}; \theta_i), y)}{\partial \mathbf{x}} \right|$ 
9:      $S_i \leftarrow S_i + S(\mathbf{x})$ 
10:  end for
11:   $S_i \leftarrow S_i / |D_i|$ 
12:  return  $(\theta_i, S_i)$ 
13: end function
14: function AGGREGATEMODELS( $\mathcal{C}, \{\theta_i\}$ )
15:   $\Theta \leftarrow \{\}$ 
16:  for each cluster  $C_k \in \mathcal{C}$  do
17:     $\theta_k \leftarrow \frac{1}{|C_k|} \sum_{i \in C_k} \theta_i$ 
18:     $\Theta \leftarrow \Theta \cup \{\theta_k\}$ 
19:  end for
20:   $\theta \leftarrow \frac{1}{|\Theta|} \sum_{\theta_k \in \Theta} \theta_k$ 
21:  return  $\theta$ 
22: end function

```

---

## A.7 THEORETICAL ANALYSIS

## A.7.1 PROOF OF PROPOSITION 1

We first state the key assumptions needed for the perturbation bounds:

**Assumptions.**

- The (normalized) graph Laplacian  $L$  is regularized (e.g.  $L \leftarrow L + \delta I$  for some  $\delta > 0$ ), or equivalently we work on the subspace orthogonal to its zero eigenvector, so that  $L$  is invertible and its Moore–Penrose inverse satisfies  $\|L^\dagger\|_2 < \infty$ .
- The perturbation magnitude is small enough that

$$\|L^\dagger\|_2 \|E_L\|_2 = \|L^\dagger\|_2 \sqrt{\epsilon_L} < 1,$$

ensuring the Neumann-series expansion for the pseudoinverse converges and the bound  $\|L'^\dagger - L^\dagger\|_2 \leq \|L^\dagger\|_2^2 \|E_L\|_2$  holds.

- The feature matrix  $X \in \mathbb{R}^{n \times f}$  has full column rank, so that its Moore–Penrose inverse  $X^\dagger$  exists and  $\|X^\dagger\|_2 < \infty$ .

**Proof.**

Saliency for a client is the (flattened) gradient of the loss w.r.t. its input features. In a first-order linearisation of a GNN we write

$$S = X^\dagger L^\dagger Y, \quad S' = X^\dagger L'^\dagger Y',$$

where

$$L' = L + E_L, \quad Y' = Y + E_Y, \quad \|E_L\|_F \leq \sqrt{\epsilon_L}, \quad \|E_Y\|_F \leq \sqrt{\epsilon_Y}.$$

Noting that  $\|E\|_2 \leq \|E\|_F$  for any matrix  $E$ , we have  $\|E_L\|_2 \leq \sqrt{\epsilon_L}$  and  $\|E_Y\|_2 \leq \sqrt{\epsilon_Y}$ .

We then bound

$$\|S' - S\|_2 = \|X^\dagger (L'^\dagger Y' - L^\dagger Y)\|_2 \leq \|X^\dagger\|_2 \left( \|L'^\dagger E_Y\|_2 + \|L'^\dagger - L^\dagger\|_2 \|Y\|_2 \right).$$

1188 We bound each term separately:

- 1189 1.  $\|L^\dagger E_Y\|_2 \leq \|L^\dagger\|_2 \|E_Y\|_2 \leq \|L^\dagger\|_2 \sqrt{\epsilon_Y}$ .
- 1190 2. By the pseudoinverse perturbation bound (valid under  $\|L^\dagger\| \sqrt{\epsilon_L} < 1$ ),
- 1191 
$$1192 \|L'^\dagger - L^\dagger\|_2 \leq \|L^\dagger\|_2^2 \|E_L\|_2 \leq \|L^\dagger\|_2^2 \sqrt{\epsilon_L}.$$

1194 Combining gives

1195 
$$1196 \|S' - S\|_2 \leq \|X^\dagger\|_2 \left( \|L^\dagger\|_2 \sqrt{\epsilon_Y} + \|L^\dagger\|_2^2 \sqrt{\epsilon_L} \|Y\|_2 \right).$$

1197 Squaring both sides yields the claimed bound:

1198 
$$1199 \|S' - S\|_2^2 \leq \|X^\dagger\|_2^2 \left( \|L^\dagger\|_2^2 \epsilon_Y + \|L^\dagger\|_2^4 \epsilon_L \|Y\|_2^2 \right).$$

1200 **Interpretation.** The first term captures sensitivity to embedding perturbations ( $\epsilon_Y$ ), while the second

1201 captures sensitivity to structural perturbations ( $\epsilon_L$ ). Both are scaled by the conditioning of the

1202 Laplacian ( $\|L^\dagger\|$ ) and the feature map ( $\|X^\dagger\|$ ), establishing that the saliency-map difference remains

1203 bounded under small perturbations.

## 1205 A.7.2 PROOF OF PROPOSITION 2

1206 We first state the additional assumptions needed for a fully rigorous bound:

### 1208 Assumptions.

- 1209 • The feature matrix  $X \in \mathbb{R}^{n \times f}$  has full column rank, so its Moore–Penrose inverse  $X^\dagger$  exists
- 1210 and  $\|X^\dagger\|_2 < \infty$ .
- 1211 • The graph Laplacian  $L$  is regularized (e.g.  $L \leftarrow L + \delta I$  with  $\delta > 0$ ) or we work on the
- 1212 subspace orthogonal to its nullspace, so that  $L^\dagger$  exists and  $\|L^\dagger\|_2 < \infty$ .
- 1213 • Perturbations are sufficiently small:

1214 
$$1215 \|X^\dagger\|_2 \|E_X\|_2 = \|X^\dagger\|_2 \sqrt{\epsilon_X} < 1, \quad \|L^\dagger\|_2 \|E_L\|_2 < 1,$$

1216 ensuring the Neumann-series expansions for pseudoinverses converge.

### 1218 Proof.

1219 Let  $M$  and  $M'$  be two GNNs differing only in their feature matrices  $X$  and  $X' = X + E_X$ , with

1220 
$$1221 \|E_X\|_F^2 \leq \epsilon_X, \quad \|E_X\|_2 \leq \sqrt{\epsilon_X}.$$

1222 The Laplacian is fixed at  $L$ , regularized so that  $L^\dagger$  is bounded. Define saliency maps

1223 
$$1224 S = X^\dagger L^\dagger Y, \quad S' = X'^\dagger L^\dagger Y',$$

1225 where  $Y$  and  $Y'$  are the corresponding node-embedding matrices, and  $\|Y' - Y\|_2 \leq \sqrt{\epsilon_Y}$ .

1226 We write

1227 
$$1228 S' - S = (X'^\dagger - X^\dagger) L^\dagger Y + X'^\dagger L^\dagger (Y' - Y).$$

1229 Using the pseudoinverse perturbation identity  $X'^\dagger - X^\dagger = -X'^\dagger E_X X^\dagger$ , we get

1230 
$$1231 (X'^\dagger - X^\dagger) L^\dagger Y = -X'^\dagger E_X X^\dagger L^\dagger Y.$$

1232 Hence by sub-multiplicativity,

1233 
$$1234 \|S' - S\|_2 \leq \|X'^\dagger\|_2 \|E_X\|_2 \|X^\dagger\|_2 \|L^\dagger\|_2 \|Y\|_2 + \|X'^\dagger\|_2 \|L^\dagger\|_2 \|Y' - Y\|_2.$$

1235 Since  $\|X'^\dagger\| \approx \|X^\dagger\|$  under  $\|X^\dagger\| \sqrt{\epsilon_X} < 1$ , and  $\|Y' - Y\| \leq \sqrt{\epsilon_Y}$ ,

1236 
$$1237 \|S' - S\|_2 \leq \|X^\dagger\|_2^2 \|L^\dagger\|_2 \|Y\|_2 \sqrt{\epsilon_X} + \|X^\dagger\|_2 \|L^\dagger\|_2 \sqrt{\epsilon_Y}.$$

1238 Squaring and absorbing the cross-term  $2\|X^\dagger\|_2^2 \|L^\dagger\|_2 \|Y\|_2 \sqrt{\epsilon_X \epsilon_Y}$  into a constant yields

1239 
$$1240 \|S' - S\|_2^2 \leq \|X^\dagger\|_2^4 \|L^\dagger\|_2^2 (\|Y\|_2^2 \epsilon_X + \epsilon_Y + \mathcal{O}(\sqrt{\epsilon_X \epsilon_Y})).$$

1241 Neglecting the higher-order cross-term gives the stated bound.

**Interpretation.** The factor  $\|X^\dagger\|_2^4$  quantifies amplification of feature perturbations,  $\|L^\dagger\|_2^2$  reflects graph-structure sensitivity,  $\|Y\|_2^2 \epsilon_X$  captures feature-to-embedding effects, and  $\epsilon_Y$  bounds embedding noise. Thus saliency-map differences remain controlled under small feature changes, validating Proposition 2.

## A.7.3 PROOF OF PROPOSITION 3: STABILITY OF SALIENCY MAPS ACROSS TASKS

**Assumptions.**

- The feature matrix  $X \in \mathbb{R}^{n \times f}$  has full column rank, so its Moore–Penrose inverse  $X^\dagger$  exists and  $\|X^\dagger\|_2 < \infty$ .
- The graph Laplacian  $L$  is regularized (e.g.  $L \leftarrow L + \delta I$  for  $\delta > 0$ ) or restricted to the subspace orthogonal to its nullspace, so  $L^\dagger$  exists and  $\|L^\dagger\|_2 < \infty$ .
- Consequently, the product  $LX$  admits a (pseudo)inverse  $(LX)^\dagger$  with  $\|(LX)^\dagger\|_2 < \infty$ .

**Proof.**

For tasks  $i$  and  $j$ , let the corresponding saliency maps be

$$S_i = X^\dagger L^\dagger Y_i, \quad S_j = X^\dagger L^\dagger Y_j,$$

where  $Y_i, Y_j$  are the node-embedding matrices. Assume  $\|Y_i - Y_j\|_2 \leq \sqrt{\epsilon_Y}$ . Then

$$S_i - S_j = X^\dagger L^\dagger (Y_i - Y_j) = (LX)^\dagger (Y_i - Y_j),$$

and by sub-multiplicativity of the spectral norm,

$$\|S_i - S_j\|_2 \leq \|(LX)^\dagger\|_2 \|Y_i - Y_j\|_2 \leq \|(LX)^\dagger\|_2 \sqrt{\epsilon_Y}.$$

Squaring both sides gives the desired bound:

$$\|S_i - S_j\|_2^2 \leq \|(LX)^\dagger\|_2^2 \epsilon_Y.$$

**Interpretation.** Here  $\|(LX)^\dagger\|_2^2$  captures how sensitive the saliency maps are to changes in the embeddings, and  $\epsilon_Y$  bounds the embedding shift between tasks. Thus, as long as task-induced embedding differences remain small, the saliency maps stay stable.